

## Session 2: MUST Correctness Checking

Dr. Matthias S. Müller (RWTH Aachen University)  
Tobias Hilbrich (Technische Universität Dresden)  
Joachim Protze (RWTH Aachen University, LLNL)

Email:

[mueller@itc.rwth-aachen.de](mailto:mueller@itc.rwth-aachen.de)  
[tobias.hilbrich@tu-dresden.de](mailto:tobias.hilbrich@tu-dresden.de)  
[protze@itc.rwth-aachen.de](mailto:protze@itc.rwth-aachen.de)

# Content

---

- Motivation
- MPI usage errors
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- Correctness tools
- MUST usage

# How many errors can you spot in this tiny example?

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);

    MPI_Recv (buf, 2, MPI_INT, size - rank, 123, MPI_COMM_WORLD, MPI_STATUS_IGNORE);

    MPI_Send (buf, 2, type, size - rank, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    return 0;
}
```

At least 8 issues in this code example

# Content

---

- Motivation
- MPI usage errors
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- Correctness tools
- MUST usage

# MPI usage errors

---

- MPI programming is error prone
- Bugs may manifest as:
  - Crashes
  - Hangs
  - Wrong results
  - Not at all! (Sleeping bugs)
- Tools help to detect these issues

# MPI usage errors (2)

---

- Complications in MPI usage:
  - Non-blocking communication
  - Persistent communication
  - Complex collectives (e.g. Alltoallw)
  - Derived datatypes
  - Non-contiguous buffers
- Error Classes include:
  - Incorrect arguments
  - Resource errors
  - Buffer usage
  - Type matching
  - Deadlocks

# Content

---

- Motivation
- MPI usage errors
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- Correctness tools
- MUST usage

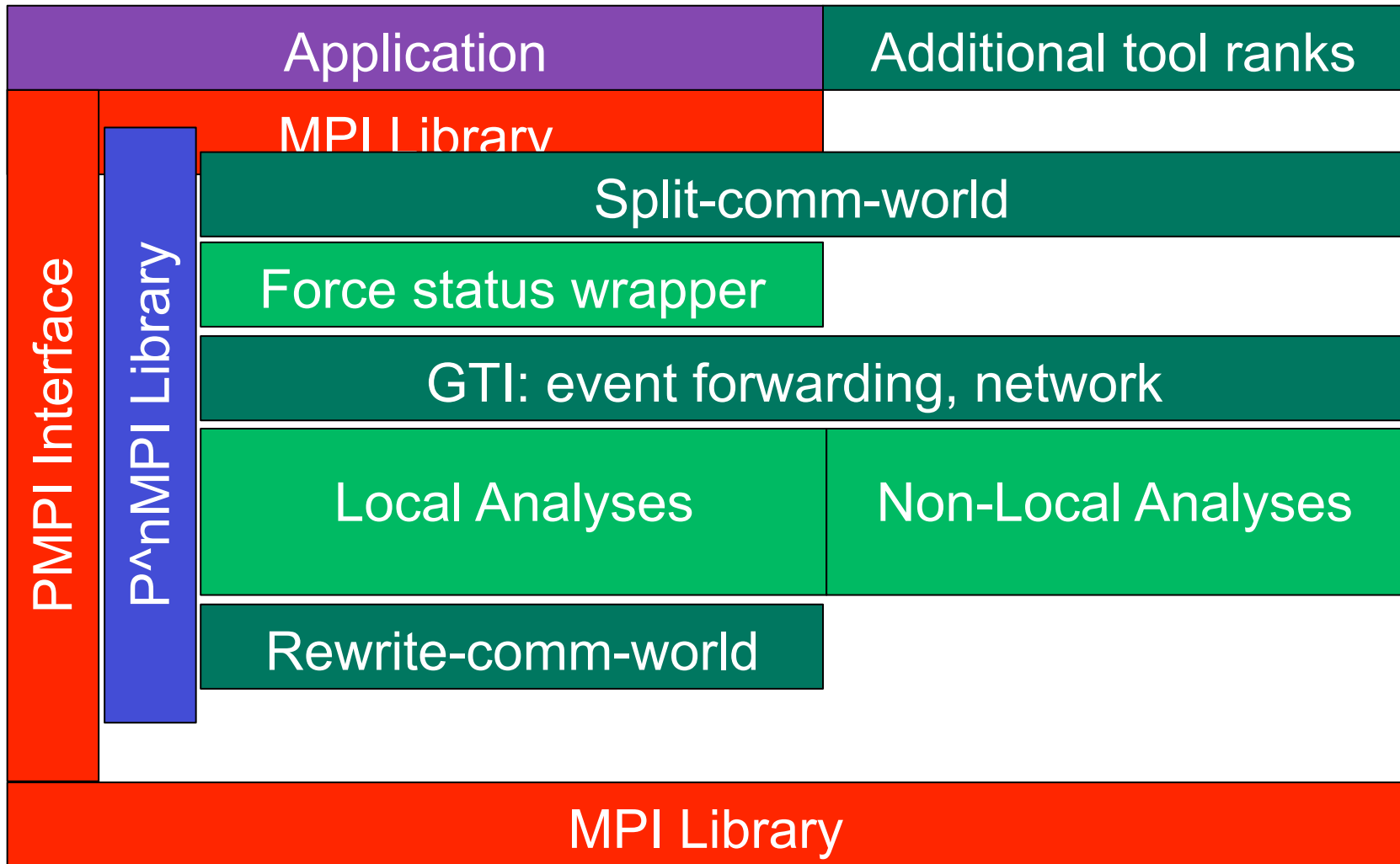
# Skipping some errors

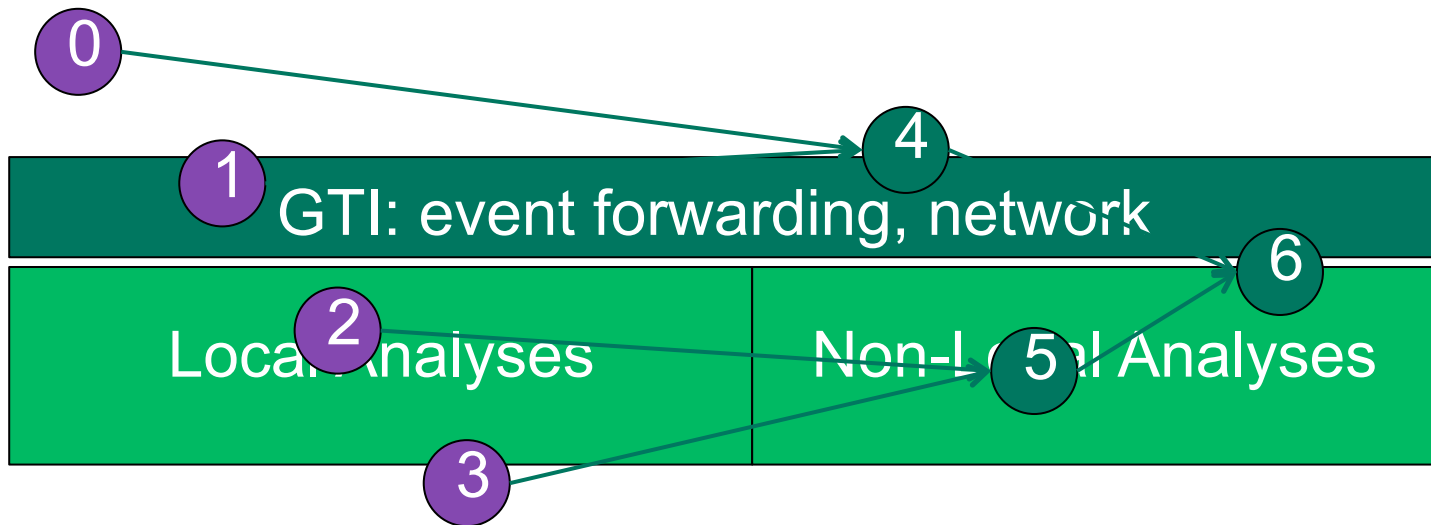
---

- Missing MPI\_Init:
  - Current release doesn't start to work, implementation in progress
- Missing MPI\_Finalize:
  - Current release doesn't terminate all analyses, work in progress
- Src/dest rank out of range (size-rank): leads to crash, use crash save version of tool



# MUST: Tool design





# Fixed some errors:

---

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);

    MPI_Recv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, MPI_STATUS_IGNORE);

    MPI_Send (buf, 2, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

# Must detects deadlocks

MUST Outputfile MUST Outputfile

file:///home/pj416018/MUST/example/MUST\_Output.html

MUST Output, starting date: Thu Nov 28 13:38:01 2013.

Rank(s)	Type	Message	From	References
	Error	The application issued a set of MPI calls that can cause a deadlock! A graphical representation of this situation is available in a <a href="#">detailed deadlock view</a> ( <a href="#">MUST_Output-files/MUST_Deadlock.html</a> ). References 1-2 list the involved calls (limited to the first 5 calls, further calls may be involved). The application still runs, if the deadlock manifested (e.g. caused a hang on this MPI implementation) you can attach to the involved ranks with a debugger or abort the application (if necessary).		References of a representative process:  reference 1 rank 0: <b>MPI_Recv</b> (1st occurrence) called from: #0 main@example.c:15  reference 2 rank 1: <b>MPI_Recv</b> (1st occurrence) called from: #0 main@example.c:15

Click for graphical representation of the detected deadlock situation.

# Graphical representation of deadlocks

MUST Outputfile MUST Outputfile  
file:///home/pj416018/MUST/example/MUST\_Output-files/MUST\_Deadlock.html

MUST Deadlock Details, date: Thu Nov 28 13:38:06 2013.

[Back to MUST error report](#)

**Message**

The application issued a set of MPI calls that can cause a deadlock! The graphs below show details on this situation. This includes a wait-for graph that shows active wait-for dependencies between the processes that cause the deadlock. Note that this process set only includes processes that cause the deadlock and no further processes. A legend details the wait-for graph components in addition, while a parallel call stack view summarizes the locations of the MPI calls that cause the deadlock. Below these graphs, a message queue graph shows active and unmatched point-to-point communications. This graph only includes operations that could have been intended to match a point-to-point operation that is relevant to the deadlock situation. Finally, a parallel call stack shows the locations of any operation in the parallel call stack. The leafs of this call stack graph show the components of the message queue graph that they span. The application still runs, if the deadlock manifested (e.g. caused a hang on this MPI implementation) you can attach to the involved ranks with a debugger or abort the application (if necessary).

**Active Communicators**

Comm: A  
MPI COMM WORLD

**Wait-for Graph**

```
graph TD; R0[0: MPI_Recv] -- "comm=A, tag=123" --> R1[1: MPI_Recv]; R1 -- "comm=" --> R0;
```

Rank 0 waits for rank 1 and vv.

**Call Stack**

```
graph TD; Main[main@example.c:15] -- "Ranks: 0-1" --> MPI[MPI_Recv];
```

Simple call stack for this example.

**Legend**

Active MPI Call (rectangle)  
Sub Operation (diamond)

A waits for B and C (solid arrows)  
A waits for B or C (dashed arrows)

**Active and Relevant Point-to-Point Messages: Overview**

**Active and Relevant Point-to-Point Messages: Callstack-view**

# Fix1: use asynchronous receive

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 2, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

Use asynchronous  
receive: (MPI\_Irecv)

# MUST detects errors in handling datatypes

MUST Outputfile			Google	
file:///home/pj416018/MUST/example/MUST_Output.html				
MUST Output, starting date: Thu Nov 28 13:50:48 2013.				
Rank(s)	Type	Message	References	References
0	Error	A receive operation uses a (datatype,count) pair that can not hold the data transferred by the send it matches! The first element of the send that did not fit into the receive operation is at (contiguous)[0](MPI_INTEGER) in the send type (consult the MUST manual for a detailed description of datatype positions). The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 2 with type:Datatype created at reference 3 is for Fortran, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)	Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix1.c:18	References of a representative process: reference 1 rank 0: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix1.c:18 reference 2 rank 1: <b>MPI_irecv</b> (1st occurrence) called from: #0 main@example-fix1.c:16 reference 3 rank 0: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix1.c:13
0-1	Error	Argument 3 (datatype) is not committed for transfer, call MPI_Type_commit before using the type for transfer! (Information on datatypeDatatype created at reference 1 is for Fortran, based on the following type(s): { MPI_INTEGER})	Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix1.c:18	References of a representative process: reference 1 rank 1: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix1.c:13
0	Error	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!  (Information on the request associated with the other communication: Request activated at reference 1) (Information on the datatype associated with the other communication: MPI_INT) The other communication overlaps with this communication at position:(MPI_INT)  (Information on the datatype associated with this communication: Datatype created at reference 2 is for Fortran, based on the following type(s): { MPI_INTEGER}) This communication overlaps with the other communication at position:(contiguous)[0](MPI_INTEGER) A graphical representation of this situation is available in a <a href="#">detailed overlap view (MUST_Output-files/MUST_Overlap_0_0.html)</a> .	Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix1.c:18	References of a representative process: reference 1 rank 0: <b>MPI_irecv</b> (1st occurrence) called from: #0 main@example-fix1.c:16 reference 2 rank 0: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix1.c:13
1	Error	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!  (Information on the request associated with the other communication: Request activated at reference 1) (Information on the datatype associated with the other communication: MPI_INT) The other communication overlaps with this communication at position:(MPI_INT)  (Information on the datatype associated with this communication: Datatype created at reference 2 is for Fortran, based on the following type(s): { MPI_INTEGER}) This communication overlaps with the other communication at position:(contiguous)[0](MPI_INTEGER) A graphical representation of this situation is available in a <a href="#">detailed overlap view (MUST_Output-files/MUST_Overlap_1_0.html)</a> .	Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix1.c:18	References of a representative process: reference 1 rank 1: <b>MPI_irecv</b> (1st occurrence) called from: #0 main@example-fix1.c:16 reference 2 rank 1: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix1.c:13
				References of a representative process:

Use of uncommitted datatype: **type**

# Fix2: use MPI\_Type\_commit

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 2, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

Commit the  
datatype before  
usage



# MUST detects errors in transfer buffer sizes / types

MUST Outputfile  
file:///home/pj416018/MUST/example/MUST\_Output.html

MUST Output, starting date: Thu Nov 28 13:51:42 2013.

Rank(s)	Type	Message	References
0	Error	<p>A receive operation uses a (datatype,count) pair that can not hold the data transferred by the send it matches! The first element of the send that did not fit into the receive operation is at (contiguous)[0](MPI_INTEGER) in the send type (consult the MUST manual for a detailed description of datatype positions). The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 2 with type:Datatype created at reference 3 is for Fortran, committed at reference 4, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)</p>	<p>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:19</p> <p>References of a representative process: reference 1 rank 0: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:19 reference 2 rank 1: <b>MPI_irecv</b> (1st occurrence) called from: #0 main@example-fix2.c:17 reference 3 rank 0: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix2.c:13 reference 4 rank 0: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix2.c:14</p>
1	Error	<p>A receive operation uses a (datatype,count) pair that can not hold the data transferred by the send it matches! The first element of the send that did not fit into the receive operation is at (contiguous)[0](MPI_INTEGER) in the send type (consult the MUST manual for a detailed description of datatype positions). The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 2 with type:Datatype created at reference 3 is for Fortran, committed at reference 4, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)</p>	<p>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:19</p> <p>References of a representative process: reference 1 rank 1: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:19 reference 2 rank 0: <b>MPI_irecv</b> (1st occurrence) called from: #0 main@example-fix2.c:17 reference 3 rank 1: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix2.c:13 reference 4 rank 1: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix2.c:14</p>
		<p>The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!</p> <p>(Information on the request associated with the other communication: Request activated at reference 1) (Information on the datatype associated with the other communication: MPI_INT)</p>	<p>Representative location: <b>MPI_Send</b> (1st occurrence)</p> <p>References of a representative process: reference 1 rank 0: <b>MPI_irecv</b> (1st occurrence) called from: #0 main@example-fix2.c:17 reference 2 rank 0:</p>

Size of sent message larger than receive buffer

# Fix3: use same message size for send and receive

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

Reduce the message  
size

# MUST detects use of wrong argument values

MUST Outputfile

file:///home/pj416018/MUST/example/MUST\_Output.html

MUST Output, starting date: Mon Dec 2 13:11:12 2013.

Rank(s)	Type	Message	References
1	Error	<p>A send and a receive operation use datatypes that do not match! Mismatch occurs at (contiguous)[0](MPI_INTEGER) in the send type and at (MPI_INT) in the receive type (consult the MUST manual for a detailed description of datatype positions). A graphical representation of this situation is available in a <a href="#">detailed type mismatch view (MUST_Output-files/MUST_Typemismatch_1.html)</a>. The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 1 with type:Datatype created at reference 3 is for Fortran, committed at reference 4, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)</p>	<p>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix3.c:19</p> <p>reference 1 rank 1: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix3.c:19</p> <p>reference 2 rank 0: <b>MPI_Recv</b> (1st occurrence) called from: #0 main@example-fix3.c:17</p> <p>reference 3 rank 1: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix3.c:13</p> <p>reference 4 rank 1: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix3.c:14</p>
0	Error	<p>A send and a receive operation use datatypes that do not match! Mismatch occurs at (contiguous)[0](MPI_INTEGER) in the send type and at (MPI_INT) in the receive type (consult the MUST manual for a detailed description of datatype positions). A graphical representation of this situation is available in a <a href="#">detailed type mismatch view (MUST_Output-files/MUST_Typemismatch_0.html)</a>. The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 1 with type:Datatype created at reference 3 is for Fortran, committed at reference 4, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)</p>	<p>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix3.c:19</p> <p>reference 1 rank 0: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix3.c:19</p> <p>reference 2 rank 1: <b>MPI_Recv</b> (1st occurrence) called from: #0 main@example-fix3.c:17</p> <p>reference 3 rank 0: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix3.c:13</p> <p>reference 4 rank 0: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix3.c:14</p>
		The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!	<p>References of a representative process:</p> <p>reference 1 rank 1:</p>

Use of Fortran type in C,  
datatype mismatch between  
sender and receiver

# Fix4: use C-datatype constants in C-code

```
#include <mpi.h>
#include <stdio.h>
```

```
int main (int argc, char** argv)
{
    int rank, size, buf[8];
```

```
    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);
```

```
    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INT, &type);
    MPI_Type_commit (&type);
```

Use the integer  
datatype intended  
for usage in C

```
    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);
```

```
    MPI_Send (buf, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);
```

```
    printf ("Hello, I am rank %d of %d.\n", rank, size);
```

```
    MPI_Finalize ();
```

```
    return 0;
```

```
}
```

# MUST detects data races in asynchronous communication

Data race between send and asynchronous receive operation

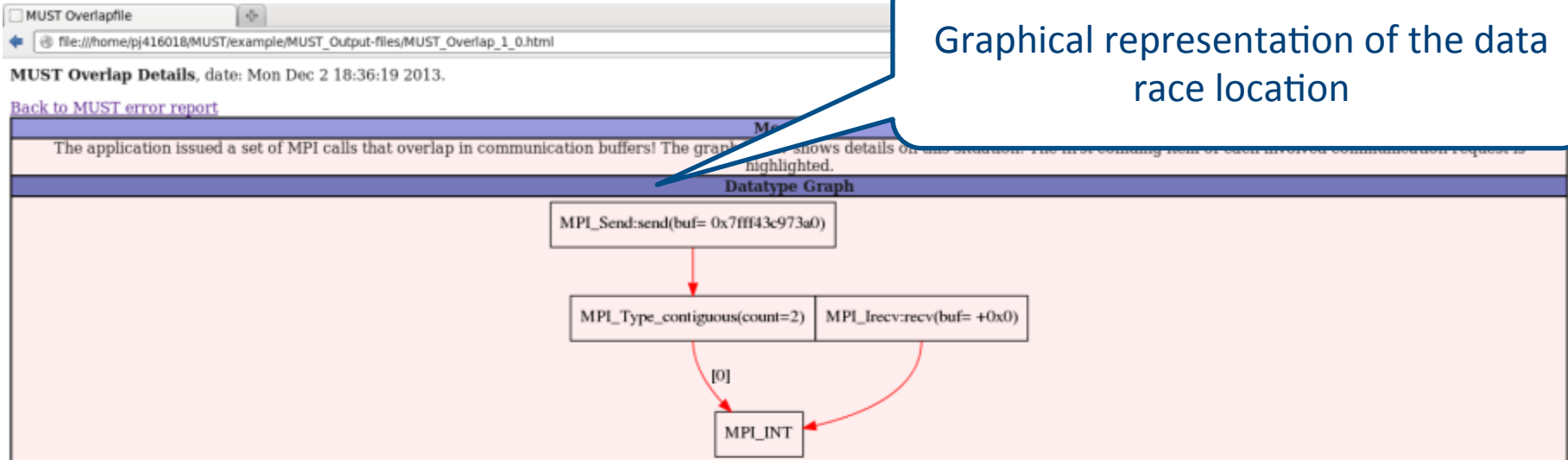
MUST Outputfile  
file:///home/pj416018/MUST/example/MUST\_Output.html

MUST Output, starting date: Mon Dec 2 18:36:19 2013.

Rank(s)	Type	Message		
1	Error	<p>The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!</p> <p>(Information on the request associated with the other communication: Request activated at reference 1)</p> <p>(Information on the datatype associated with the other communication: MPI_INT)</p> <p>The other communication overlaps with this communication at position:(MPI_INT)</p> <p>(Information on the datatype associated with this communication: Datatype created at reference 2 is for C, committed at reference 3, based on the following type(s): { MPI_INT })</p> <p>This communication overlaps with the other communication at position:(contiguous)[0](MPI_INT)</p> <p>A graphical representation of this situation is available in a <a href="#">detailed overlap view (MUST_Output-files/MUST_Overlap_1_0.html)</a>.</p>	<p>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix4.c:19</p>	<p>References of a representative process:</p> <p>reference 1 rank 1: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix4.c:17</p> <p>reference 2 rank 1: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix4.c:13</p> <p>reference 3 rank 1: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix4.c:14</p>
0-1	Error	<p>There are 1 datatypes that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these datatypes:</p> <p>-Datatype 1: Datatype created at reference 1 is for C, committed at reference 2, based on the following type(s): { MPI_INT }</p>	<p>Representative location: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix4.c:13</p>	<p>References of a representative process:</p> <p>reference 1 rank 1: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix4.c:13</p> <p>reference 2 rank 1: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix4.c:14</p>
0-1	Error	<p>There are 1 requests that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these requests:</p> <p>-Request 1: Request activated at reference 1</p>	<p>Representative location: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix4.c:17</p>	<p>References of a representative process:</p> <p>reference 1 rank 1: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix4.c:17</p>
0	Error	<p>The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!</p> <p>(Information on the request associated with the other communication: Request activated at reference 1)</p> <p>(Information on the datatype associated with the other communication: MPI_INT)</p> <p>The other communication overlaps with this communication at position:(MPI_INT)</p> <p>(Information on the datatype associated with this communication: Datatype created at reference 2 is for C, committed at reference 3, based on the following type(s): { MPI_INT })</p> <p>This communication overlaps with the other communication at position:(contiguous)[0](MPI_INT)</p> <p>A graphical representation of this situation is available in a <a href="#">detailed overlap view (MUST_Output-files/MUST_Overlap_0_0.html)</a>.</p>	<p>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix4.c:19</p>	<p>References of a representative process:</p> <p>reference 1 rank 0: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix4.c:17</p> <p>reference 2 rank 0: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix4.c:13</p> <p>reference 3 rank 0: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix4.c:14</p>

MUST has completed successfully, end date: Mon Dec 2 18:36:20 2013.

# Graphical representation of the race condition



Graphical representation of the data race location



# Fix5: use independent memory regions

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf + 4, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

Offset points to  
independent  
memory

# MUST detects leaks of user defined objects



MUST Output, starting date: Thu Nov 28 13:55:26 2013.

Rank(s)	Type	Message	From	References
0-1	Error	There are 1 datatypes that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these datatypes: -Datatype 1: Datatype created at reference 1 is for C, committed at reference 2, based on the following type(s): { MPI_INT }	Representative location: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix5.c:13	References of a representative process: reference 1 rank 0: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix5.c:13 reference 2 rank 0: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix5.c:14
0-1	Error	There are 1 requests that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these requests: -Request 1: Request activated at reference 1	Representative location: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix5.c:17	References of a representative process: reference 1 rank 0: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix5.c:17

MUST has completed successfully, end date: Thu Nov 28 13:55:26 2013.

Leak of user defined  
datatype object

- User defined objects include
  - MPI\_Comms (even by MPI\_Comm\_dup)
  - MPI\_Datatypes
  - MPI\_Groups



# Fix6: Deallocate datatype object

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INT, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf + 4, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);
    MPI_Type_free (&type);

    MPI_Finalize ();

    return 0;
}
```

Deallocate the  
created datatype

# MUST detects unfinished asynchronous communication



MUST Output, starting date: Thu Nov 28 13:55:49 2013.

Rank(s)	Type	Message	From	References
0-1	Error	<p>There are 1 requests that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these requests:</p> <p>-Request 1: Request activated at reference 1</p>	<p>Representative location: <b>MPI_irecv</b> (1st occurrence) called from: #0 main@example-fix6.c:17</p>	<p>References of a representative process: reference 1 rank 0: <b>MPI_irecv</b> (1st occurrence) called from: #0 main@example-fix6.c:17</p>

MUST has completed successfully, end date: Thu Nov 28 13:55:49 2013.

Remaining unfinished  
asynchronous receive

# Fix8: use MPI\_Wait to finish asynchronous communication

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INT, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf + 4, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    MPI_Wait (&request, MPI_STATUS_IGNORE);

    printf ("Hello, I am rank %d of %d.\n", rank, size);
    MPI_Type_free (&type);

    MPI_Finalize ();

    return 0;
}
```

Finish the  
asynchronous  
communication



**MUST Output**, starting date: Thu Nov 28 13:56:03 2013.

Rank(s)	Type	Message	From	References
	Information	MUST detected no MPI usage errors nor any suspicious behavior during this application run.		

**MUST has completed successfully**, end date: Thu Nov 28 13:56:03 2013

No further error  
detected

Hopefully this message  
applies to many  
applications

# Content

---

- Motivation
- MPI usage errors
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- Correctness tools
- MUST usage

# Classes of Correctness Tools

---

- Debuggers:

- Helpful to pinpoint any error
- Finding the root cause may be very hard
- Won't detect sleeping errors
- E.g.: gdb, TotalView, Alinea DDT

- Static Analysis:

- Compilers and Source analyzers
- Typically: type and expression errors
- E.g.: MPI-Check

- Model checking:

- Requires a model of your applications
- State explosion possible
- E.g.: MPI-Spin

# Strategies of Correctness Tools

---

- Runtime error detection:

- Inspect MPI calls at runtime
- Limited to the timely interleaving that is observed
- Causes overhead during application run
- E.g.: Intel Trace Analyzer, Umpire, Marmot, MUST

- Formal verification:

- Extension of runtime error detection
- Explores ALL possible timely interleavings
- Can detect potential deadlocks or type mismatches that would otherwise not occur in the presence of a tool
- For non-deterministic applications exponential exploration space
- E.g.: ISP

# Content

---

- Motivation
- MPI usage errors
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- Correctness tools
- MUST usage



# MUST Usage

---

- 1) Compile and link application as usual
    - Link against the shared version of the MPI lib (Usually default)
  - 2) Replace “mpiexec” with “mustrun”
    - E.g.: *mustrun -np 4 myApp.exe input.txt output.txt*
  - 3) Inspect “MUST\_Output.html” in run directory
    - “MUST\_Output/MUST\_Deadlock.dot” exists in case of deadlock
    - Visualize with: *dot -Tps MUST\_Deadlock.dot -o deadlock.ps*
- The mustrun script will use an extra process for non-local checks (Invisible to application)
  - I.e.: “mustrun -np 4 ...” will issue a “mpirun -np 5 ...”
  - Make sure to allocate the extra task in batch jobs

- Local checks:

- Integer validation
- Integrity checks (pointers valid, etc.)
- Operation, Request, Communicator, Datatype, Group usage
- Resource leak detection
- Memory overlap checks

- Non-local checks:

- Collective verification
- Lost message detection
- Type matching (For P2P and collectives)
- Deadlock detection (with root cause visualization)

# MUST - Features: Scalability

---

- Local checks largely scalable
- Non-local checks:
  - Current default uses a central process
    - This process is an MPI task taken from the application
    - Limited scalability ~100 tasks (Depending on application)
  - Distributed analysis available (tested with 10k tasks)
    - Uses more extra tasks (10%-100%)
- Recommended: Logging to an HTML file
- Uses a scalable tool infrastructure
- Tool configuration happens at execution time